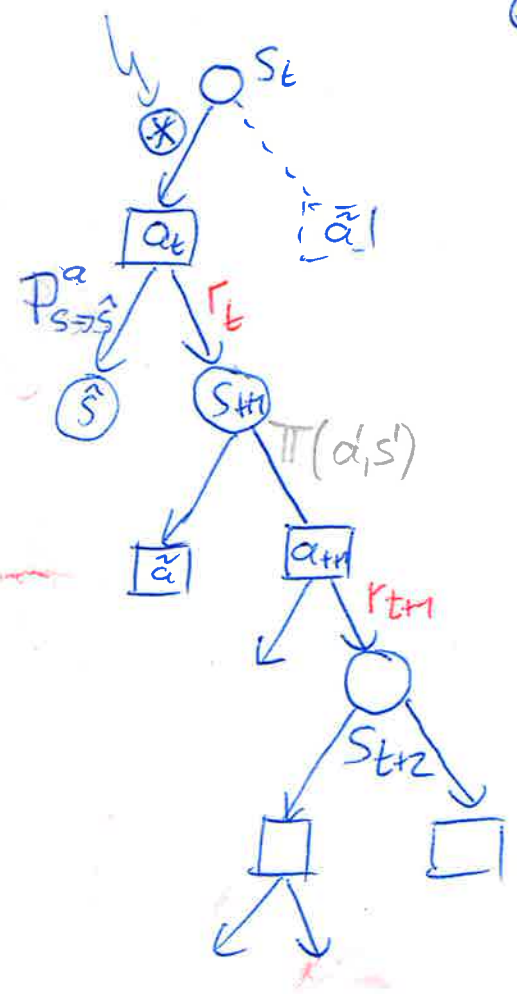


we start here



total reward collected in single trial starting in \$s\$ with action \$a\_t\$

$$\begin{aligned}
 R^{tot}(s_t, a_t) &= r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \gamma^3 r_{t+3} + \dots \\
 &= r_t + \gamma [r_{t+1} + \gamma r_{t+2} + \gamma^2 r_{t+3} + \dots] \\
 &= r_t + \gamma \cdot R^{tot}(s_{t+1}, a_{t+1})
 \end{aligned}$$

total reward (single trial) starting from \$s' = s\_{t+1}\$ with \$a\_{t+1}\$

now we look at diagram to calculate expectation

$$\begin{aligned}
 E(R^{tot}(s_t, a_t)) &= E(r_t + \gamma R^{tot}(s_{t+1}, a_{t+1})) \\
 &= \sum_{s'} P_{s \to s'}^{a_t} [R_{s \to s'}^{a_t} + \gamma E(R^{tot}(s'))] \\
 &= \sum_{s'} P_{s \to s'}^{a_t} [R_{s \to s'}^{a_t} + \gamma \sum_{a'} \Pi(a', s') E(R^{tot}(s', a'))]
 \end{aligned}$$

↑ starting in \$s'\$

$$Q(s_t, a_t) = \sum_{s'} P_{s \to s'}^{a_t} [R_{s \to s'}^{a_t} + \gamma \sum_{a'} \Pi(a', s') Q(s', a')]$$

all future states and actions, given \$(s\_t, a\_t)\$

# Blackboard RL-5 - SARSA

5

from diagram

$$Q(s, a) \approx r_t + \underset{\substack{\text{discount} \\ \downarrow}}{\gamma} \cdot Q(s', a')$$

$$0 \approx r_t + \gamma \cdot Q(s', a') - Q(s, a)$$

proposed update

$$\Delta Q(s, a) = \eta [r_t + \gamma \cdot Q(s', a') - Q(s, a)]$$

check:

$$\text{if } r_t > \underbrace{\gamma \cdot Q(s', a') - Q(s, a)}_{\substack{\text{expected reward} \\ \text{for this transition}}} \Rightarrow \text{increase } Q(s, a)$$

↑  
actual  
reward

# Blackboard 6A : Consistency of Q-values

(6A)

assumption:

- ~~the~~ set of Q-values is fixed ("frozen")

$$\underline{0} = E[\Delta Q(s, a) | s, a]$$

↑ expectation of next possible states/actions  
[with frozen policy consistent with frozen Q-values]

⇒ Q-values consistent with Bellman eq.

proof : expectation over all possible paths starting in  $(s, a)$

$$0 = \frac{1}{\eta} E[\Delta Q(s, a)] =$$

$$= E\left[\frac{1}{\eta} \Delta Q(s, a)\right]$$

$$= E\left[r_t + \gamma Q(s', a') - Q(s, a)\right]$$

$$= E\left[r_t + \gamma \underset{\substack{\uparrow \\ \text{frozen}}}{Q(s', a')}\right] - \underset{\substack{\uparrow \\ \text{frozen!} \\ \text{trivial expectation}}}{Q(s, a)}$$

shift  $Q(s, a)$  to left-hand side

$$\underline{Q(s, a)} = \sum_{s'} P_{s \rightarrow s'}^a \left[ R_{s \rightarrow s'}^a + \sum_{a'} \pi(s', a') Q(s', a') \right]$$

⇒ Bellman equation. ◻

# Consistency of fluctuating online SARSA with Bellman equation

hypothesis:  $0 = \langle \Delta \hat{Q}(s,a) \rangle = \eta \langle r_t + \gamma \hat{Q}(s',a') - \hat{Q}(s,a) \rangle$

*cut* (under  $r_t + \gamma \hat{Q}(s',a')$ )

*shift* (under  $\hat{Q}(s,a)$ )

$\langle \hat{Q}(s,a) \rangle = \langle r_t + \gamma \hat{Q}(s',a') \rangle$

$\uparrow$  fluctuates

$\uparrow$  fluctuates

temporal average over many "trials" ( $N \rightarrow \infty$ )

$\langle \hat{Q}(s,a) \rangle = \sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma \sum_{a'} \langle \pi(s',a') \cdot \hat{Q}(s',a') \rangle]$

*Problem:*  $\pi^a$  depends on  $Q$

Solution: ① if  $\eta$  is small, the fluctuations of  $\hat{Q}$  are small and fluctuations of policy  $\pi^a$  are "even smaller"

② consider  $\pi^a$  fixed for small enough  $\eta$   
 $\Rightarrow$  move  $\pi$  out:  $\langle \pi^a \hat{Q} \rangle \sim \pi^a \langle Q \rangle$

$\langle \hat{Q}(s,a) \rangle = \sum_{s'} P_{s \rightarrow s'}^a [R_{s \rightarrow s'}^a + \gamma \sum_{a'} \pi^a(s',a') \langle \hat{Q}(s',a') \rangle]$

$\downarrow$

$\langle \hat{Q}(s,a) \rangle = Q(s,a)$  solves Bellman equation  $Q(s',a') = R(s',a')$

## remarks

① \*example of  $\pi^a$  "even smaller":  $\epsilon$ -greedy  
 only rank-order of  $Q$ -values matters: best / 2<sup>nd</sup> best / ...  
 if fluctuations  $|\Delta \hat{Q}(s',a')| \ll Q(s',\text{best}) - Q(s',\text{2<sup>nd</sup> best})$   
 then  $\pi^a$  remains stable!

② evaluation of averages - look at graph  
 - if in state  $s'$  all remaining averages are "given  $s'$ "